

Machine Learning

Haoyu Zhen

March 4, 2022

Contents

1	Foundations	2
1.1	Model evaluation:	2
1.2	Performance	2
1.3	Bias-Variance Decomposition	2
2	Regression	2
2.1	Linear Regression	2
2.2	Ridge Regression	3
2.3	Lasso Regression	4
2.4	Logistic Regression	4

Acknowledgement

These Notes contain material developed and copyright by:

- *AI-2611, Machine Learning*, © 2022 by Bingbing Ni, Shanghai Jiao Tong University.
- *Understanding Machine Learning*, © 2014 by Shai Shalev-Shwartz and Shai Ben-David.
- *The Elements of Statistical Learning*, © 2017 by Trevor Hastie, Robert Tibshirani and Jerome Friedman.

1 Foundations

1.1 Model evaluation:

Hold-out, cross validation and bootstrap.

For cross validation, we often let the numbers of the folds be 10. And in bootstrap, the equation $\lim_{n \rightarrow \infty} (1 - 1/m)^m = 1/e$ is used to analyse the probability.

1.2 Performance

Now we consider that:

	prediction+	prediction-
Actual	1	0
1	TP	FP
0	FN	TN

Definition 1.1 (Sensitivity and FPR).

$$TPR = \frac{TP}{TP + FN}, FPR = \frac{FP}{TN + FP}.$$

Then we introduce **ROC** space and **AUC**.

Definition 1.2 (Precision and recall).

$$precision = \frac{TP}{TP + FP}, recall = \frac{TP}{TP + FN}.$$

$$F_\beta = \frac{(1 + \beta^2) \times P \times R}{\beta^2 \times P + R}.$$

β depends on the preference of Precision and Recall.

1.3 Bias-Variance Decomposition

$$\begin{aligned} E(f; D) &= bias^2(x) + var(x) + \varepsilon^2 \\ &= (\bar{f}(x) - y)^2 + \mathbb{E}_D[f(x; D) - \bar{f}(x)] + \mathbb{E}_D[(y_D - y)^2] \end{aligned}$$

2 Regression

2.1 Linear Regression

The hypothesis class of linear regression predictors is simply the set of linear functions,

$$\mathcal{H}_{reg} = \{ \mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \}.$$

Intuitively,

$$\mathcal{L}_S(h) = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}) - \mathbf{y})^2, \forall h \in \mathcal{H}_{reg}.$$

To minimize the loss function, we need to solve $A\mathbf{w} = \mathbf{b}$ where $A \stackrel{\text{def}}{=} \sum \mathbf{x}_i \mathbf{x}_i^T = XX^T$ and $\mathbf{b} \stackrel{\text{def}}{=} \sum y_i \mathbf{x}_i = X^T \mathbf{y}$. If A is invertible then the solution is $w = A^{-1}\mathbf{b}$.

Theorem 2.1.

$$\omega = (X^T X)^{-1} X^T \mathbf{y}.$$

If the training instances do not span the entire space of \mathbb{R}^d then A is not invertible.

Theorem 2.2. Using A 's eigenvalue decomposition, we could write A as VD^+V^T where D is a diagonal matrix and V is an orthonormal matrix. Define D^+ to be the diagonal matrix such that $D_{i,i}^+ = 0$ if $D_{i,i} = 0$ otherwise $D_{i,i}^+ = 1/D_{i,i}$. Then,

$$A\hat{\mathbf{w}} = \mathbf{b}$$

where $\hat{\mathbf{w}} = VD^+V^T\mathbf{b}$

Proof.

$$A\hat{\mathbf{w}} = AA^+\mathbf{b} = VDV^TVD^+V^T\mathbf{b} = VDD^+V^T\mathbf{b} = \sum_{i:D_{i,i} \neq 0} \mathbf{v}_i \mathbf{v}_i^T \mathbf{b}.$$

That is, $A\hat{\mathbf{w}}$ is the projection of \mathbf{b} onto the span of those vectors \mathbf{v}_i for which $D_{i,i} \neq 0$. Since the linear span of $\mathbf{x}_1, \dots, \mathbf{x}_m$ is the same as the linear span of those \mathbf{v}_i , and \mathbf{b} is in the linear span of the \mathbf{x}_i , we obtain that $A\hat{\mathbf{w}} = \mathbf{b}$, which concludes our argument. \square

Remark 2.1. Indeed we always use the **Gradient Descent** method to optimize the loss function.

Linear regression for polynomial regression tasks $\mathcal{H}_{poly}^n = \{x \mapsto p(x)\}$ where $\psi(x) = (1, x, x^2, \dots, x^n)$ and $p(\psi(x)) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$.

2.2 Ridge Regression

To ameliorate the effect of the invertible matrix, we could introduce the regularization.

Definition 2.1 (Regularized Loss).

$$R(w) = \lambda \|w\|^2.$$

Now the loss function reads:

$$\mathcal{L} = \mathcal{L}_S(w) + R(w) = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}) - \mathbf{y})^2 + \lambda \|w\|^2.$$

2.3 Lasso Regression

2.4 Logistic Regression